

# *In silico* prediction of Genomic Islands in microbial genomes

Georgios S. Vernikos  
Selwyn College  
University of Cambridge

April 2008

A dissertation submitted for  
the degree of Doctor of Philosophy  
at the University of Cambridge



## **Declaration**

The work presented in this dissertation was carried out at the Wellcome Trust Sanger Institute in Cambridge between April 2005 and April 2008. This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration except where specifically indicated in the text. No part of this dissertation has been, or is being submitted for a degree or diploma or other qualification at this or any other university.

Georgios S. Vernikos  
Cambridge  
April 2008



## Summary

Large inserts of horizontally acquired DNA that contain functionally related genes with limited phylogenetic distribution are often referred to as Genomic Islands (GIs). Integration of GIs can in a single step transform a microbial organism, changing radically the way it interacts with its niche, a process that is often referred to as “evolution in quantum leaps”. A key aspect in the prediction of GIs is that at the time of their integration in the new host chromosome, their composition reflects mainly the composition of the donor, rather than the host.

The first part of this thesis concerns the design, development and implementation of a novel algorithm, that of the Interpolated Variable Order Motifs, for the composition-based prediction of GIs. This algorithm exploits compositional biases using variable order motif distributions and captures more reliably the local composition of a sequence compared with fixed order methods, overcoming the limitations of the latter. Furthermore, for the optimal localization of the boundaries of each predicted region, the Hidden Markov Model theory was implemented in a change point detection framework, predicting more accurately the true insertion point of candidate GIs.

In the second part of this thesis whole genome based comparative and phylogenetic techniques were used to study the acquisition of horizontally acquired genes in the *Salmonella* lineage in a time dependent manner. The compositional amelioration process was modelled and the relative time of acquisition of those genes was determined on different branches of the *S. enterica* phylogenetic tree.

The aim of the third part of this thesis is to explicitly quantify and model the contribution of genomic features to the GI structure, under a probabilistic framework. A hypothesis free, bottom-up search was implemented and identified approximately 700 genomic regions, including

both GIs and randomly sampled regions from three different genera that form my training dataset. A Machine Learning approach was used to exploit the above dataset and study the structural variation of GIs.

The last part of this thesis focuses on the experimental validation of the *in silico* predictions made on a newly sequenced bacterial genome. Applying a PCR-based protocol, the presence and absence of the predicted candidate islands was probed in seventeen unsequenced closely and distantly related strains. The true borders of the predicted islands were confirmed by sequencing across the boundary site in strains lacking the island.



## Acknowledgements

Firstly, my deepest personal thanks to my supervisor Julian Parkhill for being the busiest and at the same time the most available “teacher” that shaped, in a very profound and enjoyable way, my academic career; Julian has a remarkable ability to inspire people working next to him.

I would also like to thank members of my thesis committee, George Salmond for discussion on various biological aspects of my project, Richard Durbin and Alex Bateman for ideas and discussion on algorithmic techniques and Thomas Down for replying to my emails using pseudo-code! – It made our communication much easier and fun.

Many thanks to David Carter for making the source code of the biojava implementation of the Relevance Vector Machine available, Nick Thomson for long, late night philosophical (and not only) discussions on microbial genomes; Stephen Bentley for biological discussion and for contributing to a key step of my future career path, Matthew Holden for discussion on various aspects of microbial evolutionary dynamics and the other members of team 81 for advice on different steps of this project.

I also thank: Helena Seth-Smith and Paul Scott, for being extremely patient with my ignorance on lab-based techniques and Andrew Jackson for discussion on aspects of the phylogenetic analysis; Xavier Didelot for advice on whole genome sequence alignments and Matthew Avison for providing the *Stenotrophomonas maltophilia* strains; Manolis Dermitzakis for my very first, off-the-record, coffee-interview that took place in Greece; the Wellcome Trust Sanger Institute for my PhD studentship.

Special thanks to my parents, Stelios and Maria, for letting me take apart their washing machine, in order for me to discover what lies beneath; their very liberal approach to my intellectual curiosity shaped in a very profound way my academic and non-academic way of thinking; my brother, Christos, for convincing me that different is nice.

My love to Andriana (my wife to be) for standing by me no matter what, without whom at the end of this 4-year academic journey, I would only be a scientist.





## Contents

<b>Declaration</b>	<b>iii</b>
<b>Summary</b>	<b>v</b>
<b>Acknowledgements</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Horizontal Gene Transfer .....	1
1.2 Genomic Islands .....	9
1.2.1 SaPIs .....	14
1.2.2 LEE .....	16
1.2.3 SPI-7 .....	19
1.2.4 SGI-1 .....	22
1.2.5 cag PAI .....	24
1.2.6 Symbiosis island .....	26
1.2.7 SXT .....	28
1.2.8 HPI .....	30
1.2.9 SPI-1, 2, 3, 4 .....	34
1.2.10 Metabolic Island .....	40
1.3 <i>In silico</i> prediction of GIs .....	41
<b>2 Alien_Hunter algorithm</b>	<b>48</b>
2.1 Introduction .....	48
2.2 Methods .....	50
2.2.1 Interpolated Variable Order Motifs .....	50
2.2.2 Relative entropy .....	53
2.2.3 Score threshold .....	55
2.2.4 Change-point detection .....	59
2.2.5 Reciprocal FASTA .....	64
2.3 Results .....	66
2.3.1 Manually curated HGT dataset .....	66
2.3.2 Three novel SPIs .....	67
2.3.3 Predicted boundary optimization .....	72
2.3.4 Performance benchmarking .....	74
2.4 Discussion .....	77
<b>3 Genetic flux over time</b>	<b>81</b>
3.1 Introduction .....	81
3.2 Methods .....	83
3.2.1 Whole Genome Alignment .....	83
3.2.2 Phylogenetic tree building methods .....	88
3.2.2.1 UPGMA .....	89
3.2.2.2 Maximum Parsimony .....	90
3.2.2.3 Bayesian inference .....	91
3.2.2.4 Neighbor – Joining .....	92
3.2.2.5 Maximum Likelihood .....	93

---

3.2.3	Nucleotide substitution models	97
3.2.3.1	Jukes-Cantor model	97
3.2.3.2	Kimura – 2 parameter model	101
3.2.3.3	F84 model	102
3.2.3.4	Substitution rate variation	102
3.2.3.5	Parameter estimation	105
3.2.4	Relative time of HGT events	107
3.2.5	Compositional analysis	111
3.2.6	Orthologous genes	111
3.3	Results	112
3.3.1	Time distribution of PHA genes	112
3.3.2	Functional analysis of PHA genes	115
3.3.3	Compositional analysis	118
3.4	Discussion	127
3.5	Conclusions	132
<b>4</b>	<b>Resolving the structure of Genomic Islands</b>	<b>134</b>
4.1	Introduction	134
4.2	Methods	138
4.2.1	Genomic Dataset	139
4.2.2	Best reciprocal FASTA	139
4.2.3	Multiple Sequence Alignments	142
4.2.4	Phylogenetic analysis	142
4.2.5	Comparative analysis	144
4.2.6	Random sampling	155
4.2.7	Structural annotation	155
4.2.8	Machine Learning	157
4.2.9	ROC curve	162
4.2.10	Cross-Validation	163
4.3	Results	163
4.3.1	GI structural models	164
4.3.1.1	Genus-specific	165
4.3.1.2	Cross-genus	171
4.3.2	Prediction accuracy	175
4.4	Discussion	177
<b>5</b>	<b>Experimental validation of the predictions</b>	<b>189</b>
5.1	Introduction	189
5.2	Methods	190
5.2.1	<i>In silico</i> prediction of GIs	191
5.2.2	Comparative analysis	191
5.2.3	Principle of the experimental approach	192
5.2.4	DNA purification	193
5.2.5	Primer design	195
5.2.6	Polymerase Chain Reaction – PCR	196
5.2.7	Gel electrophoresis	197
5.2.8	Sequencing	198
5.3	Results	198

5.3.1	Genomic Island candidates.....	198
5.3.1.1	Genomic Island 1 .....	198
5.3.1.2	Genomic Island 16.....	201
5.3.1.3	Genomic Island 4 .....	203
5.3.1.4	Genomic Island 12.....	205
5.3.1.5	Genomic Island 14.....	207
5.3.1.6	Genomic Island 15.....	210
5.3.1.7	Genomic Island 20.....	214
5.3.1.8	Genomic Island 7.....	216
5.3.2	Performance benchmarking .....	217
5.3.2.1	Prediction accuracy .....	217
5.3.2.2	Boundary accuracy .....	218
5.4	Discussion.....	219
<b>6</b>	<b>Discussion</b>	<b>224</b>
6.1	Conclusions .....	224
6.2	Future work .....	227
6.3	Final remarks .....	232
	<b>Bibliography</b>	<b>234</b>